

Big Data Infrastructure

CS 489/698 Big Data Infrastructure (Winter 2016)

Week 10: Mutable State (2/2)

March 17, 2016

Jimmy Lin

David R. Cheriton School of Computer Science

University of Waterloo

These slides are available at <http://lintool.github.io/bigdata-2016w/>

This work is licensed under a Creative Commons Attribution-Noncommercial-Share Alike 3.0 United States
See <http://creativecommons.org/licenses/by-nc-sa/3.0/us/> for details



The Fundamental Problem

- We want to keep track of *mutable* state in a *scalable* manner
- Assumptions:
 - State organized in terms of many “records”
 - State unlikely to fit on single machine, must be distributed

(note: much of this material belongs in a distributed systems or databases course)

Motivating Scenarios

- Money shouldn't be created or destroyed:
 - Alice transfers \$100 to Bob and \$50 to Carol
 - The total amount of money after the transfer should be the same
- Phantom shopping cart:
 - Bob removes an item from his shopping cart...
 - Item still remains in the shopping cart
 - Bob refreshes the page a couple of times... item finally gone

Motivating Scenarios

- People you don't want seeing your pictures:
 - Alice removes mom from list of people who can view photos
 - Alice posts embarrassing pictures from Spring Break
 - Can mom see Alice's photo?
- Why am I still getting messages?
 - Bob unsubscribes from mailing list
 - Message sent to mailing list right after
 - Does Bob receive the message?

Three Core Ideas

- Partitioning (sharding)

- For scalability
- For latency

Need distributed transactions!

- Replication

- For robustness (availability)
- For throughput

Need replica coherence protocol!

- Caching

- For latency

Need cache coherence protocol!

How to address?



Relational Databases

... to the rescue!

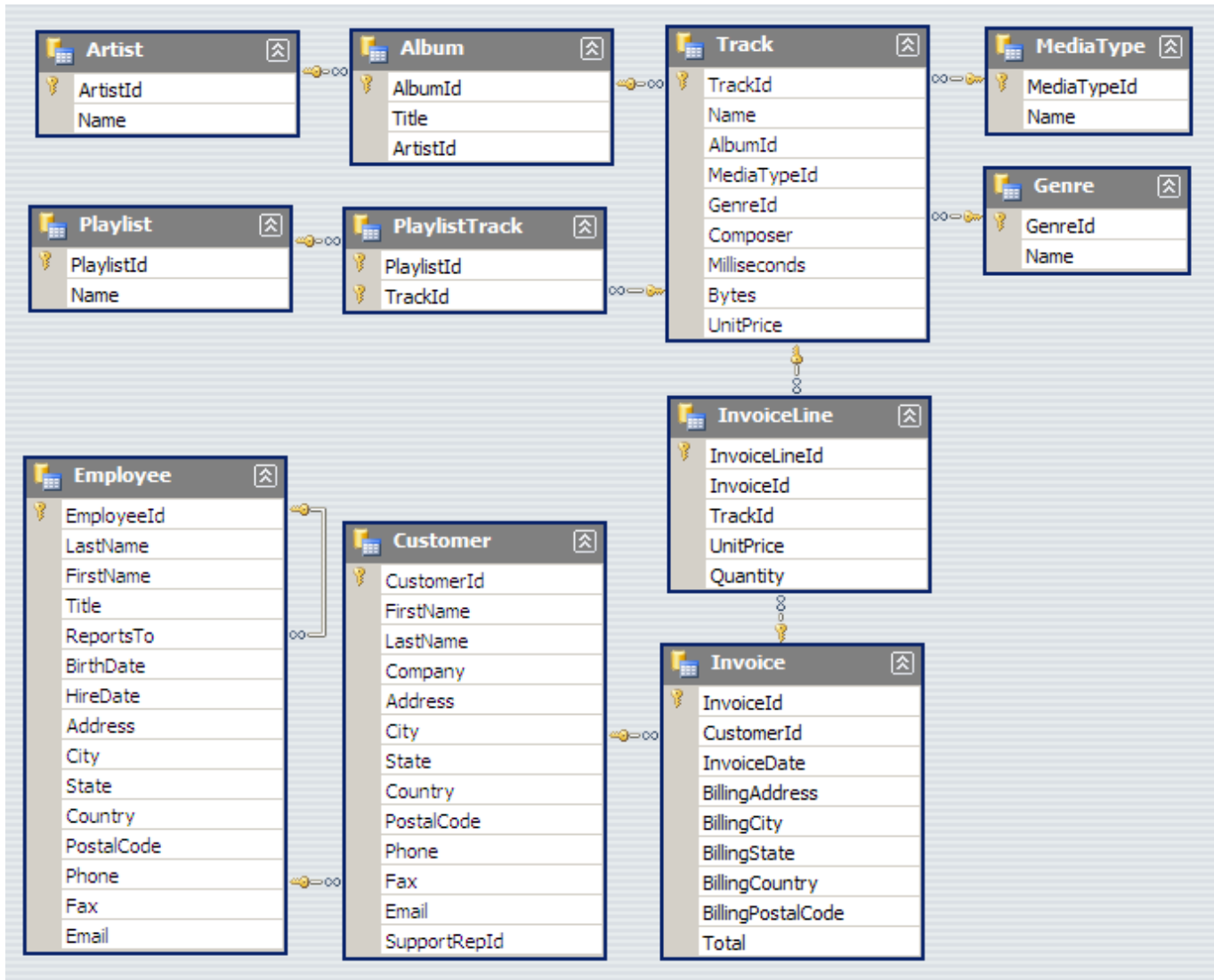
What do RDBMSes provide?

- Relational model with schemas
- Powerful, flexible query language
- Transactional semantics: ACID
- Rich ecosystem, lots of tool support



RDBMSes: Pain Points

#1: Must design up front, painful to evolve



Note: Flexible design doesn't mean *no* design!

#2: Pay for ACID!



#3: Cost!



What do RDBMSes provide?

- Relational model with schemas
- Powerful, flexible query language
- Transactional semantics: ACID
- Rich ecosystem, lots of tool support

What if we want *a la carte*?



Features *a la carte*?

- What if I'm willing to give up consistency for scalability?
- What if I'm willing to give up the relational model for something more flexible?
- What if I just want a cheaper solution?

Three Core Ideas

○ Partitioning (sharding)

- For scalability
- For latency

Need distributed transactions!

○ Replication

- For robustness (availability)
- For throughput

Need replica coherence protocol!

○ Caching

- For latency

Need cache coherence protocol!

Motivating application?

How do RDBMSes do it?

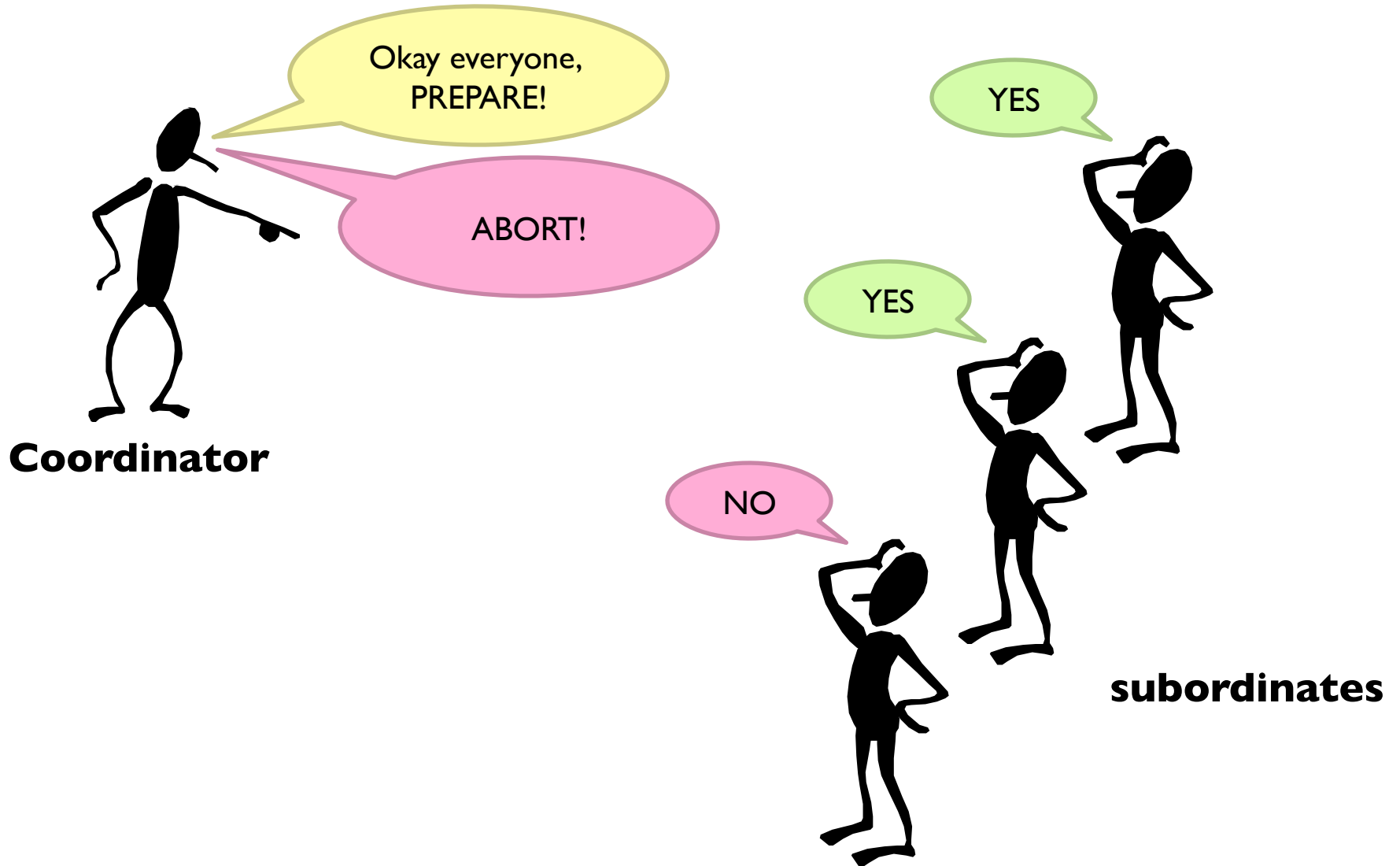
- Transactions on a single machine: (relatively) easy!
- Partition tables to keep transactions on a single machine
 - Example: partition by user
- What about transactions that require multiple machine?
 - Example: transactions involving multiple users

Solution: Two-Phase Commit

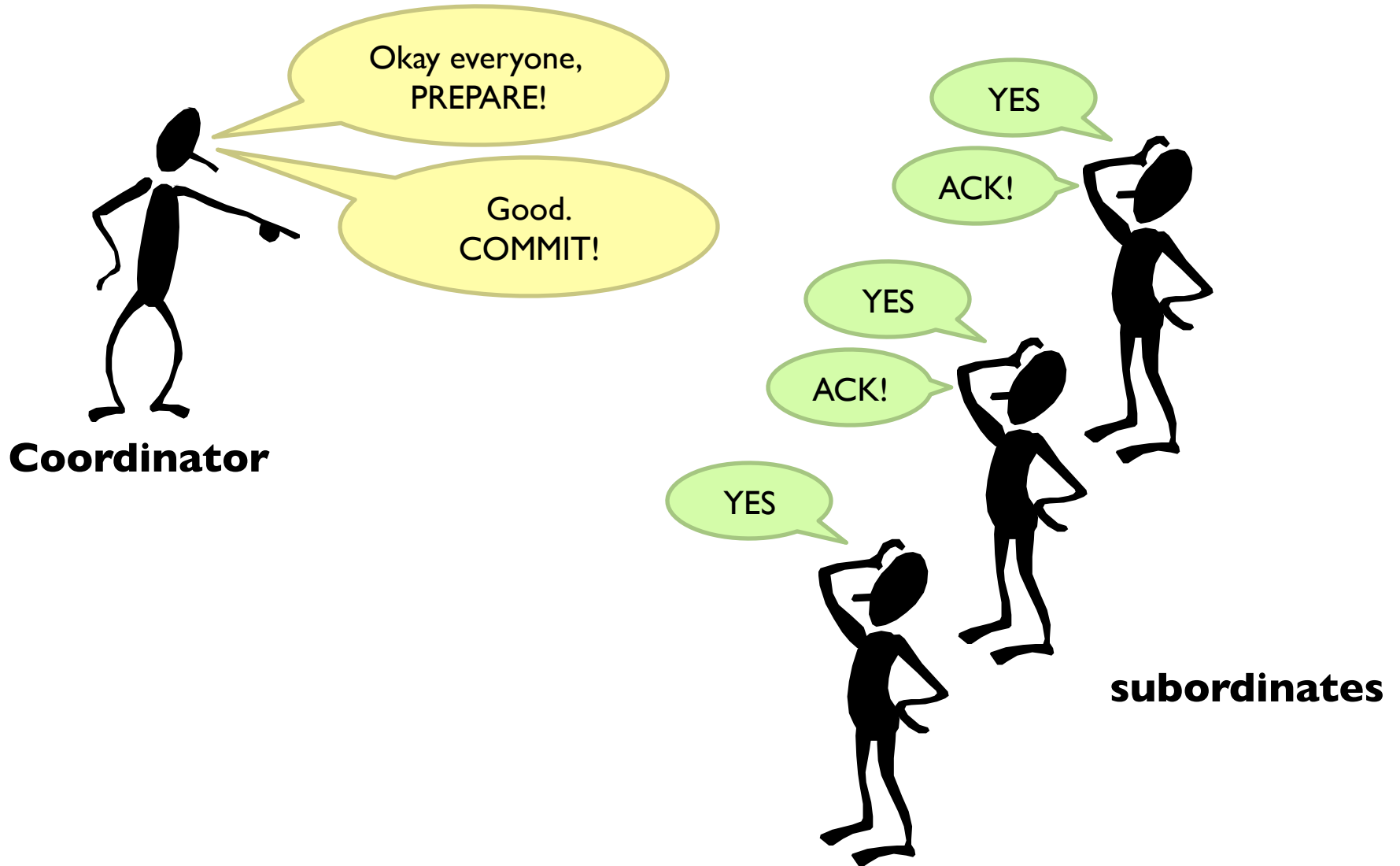
2PC: Sketch



2PC: Sketch



2PC: Sketch



2PC: Assumptions and Limitations

- Assumptions:

- Persistent storage and write-ahead log at every node
- WAL is never permanently lost

- Limitations:

- It's blocking and slow
- What if the coordinator dies?

Beyond 2PC: Paxos!
(details beyond scope of this course)

Remember this!

Key-Value Stores: Operations

- Very simple API:
 - Get – fetch value associated with key
 - Put – set value associated with key
- Optional operations:
 - Multi-get
 - Multi-put
 - Range queries
- Consistency model:
 - Atomic puts (usually)
 - Cross-key operations: who knows?

“Unit of Consistency”

- Single record:
 - Relatively straightforward
 - Complex application logic to handle multi-record transactions
- Arbitrary transactions:
 - Requires 2PC
- Middle ground: entity groups
 - Groups of entities that share affinity
 - Co-locate entity groups
 - Provide transaction support within entity groups
 - Example: user + user’s photos + user’s posts etc.

Where have we learned this trick before?

Three Core Ideas

- Partitioning (sharding)

- For scalability
- For latency

Need distributed transactions!

- Replication

- For robustness (availability)
- For throughput

Need replica coherence protocol!

- Caching

- For latency

Need cache coherence protocol!

CAP “Theorem” (Brewer, 2000)

Consistency

Availability

Partition tolerance

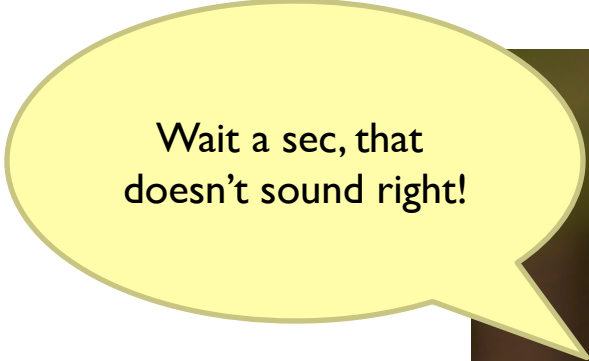
... pick two

CAP Tradeoffs

- CA = consistency + availability
 - E.g., parallel databases that use 2PC
- AP = availability + tolerance to partitions
 - E.g., DNS, web caching

Is this helpful?

- CAP not really even a “theorem” because vague definitions
 - More precise formulation came a few years later



Wait a sec, that
doesn't sound right!



Abadi Says...

- CP makes no sense!
- CAP says, in the presence of P, choose A or C
 - But you'd want to make this tradeoff even when there is no P
- Fundamental tradeoff is between consistency and latency
 - Not available = (very) long latency

Replication possibilities

- Update sent to all replicas at the same time
 - To guarantee consistency you need something like Paxos
- Update sent to a master
 - Replication is synchronous
 - Replication is asynchronous
 - Combination of both
- Update sent to an arbitrary replica

All these possibilities involve tradeoffs!
“eventual consistency”

Move over, CAP

- PACELC (“pass-elk”)
- PAC
 - If there’s a partition, do we choose A or C?
- ELC
 - Otherwise, do we choose latency or consistency?

To: All Graduate Students

Due to a recent incident, we would like to remind all Grad Students that refreshments provided in communal areas during an event are for attendees of that event only.

Please vacate the communal area and do not consume the refreshments unless you have been specifically invited to participate.

To avoid any misunderstanding, you are only invited if you received a specific invitation by e-mail or if it was arranged by your supervisor for you to attend.

Thank you for your cooperation,

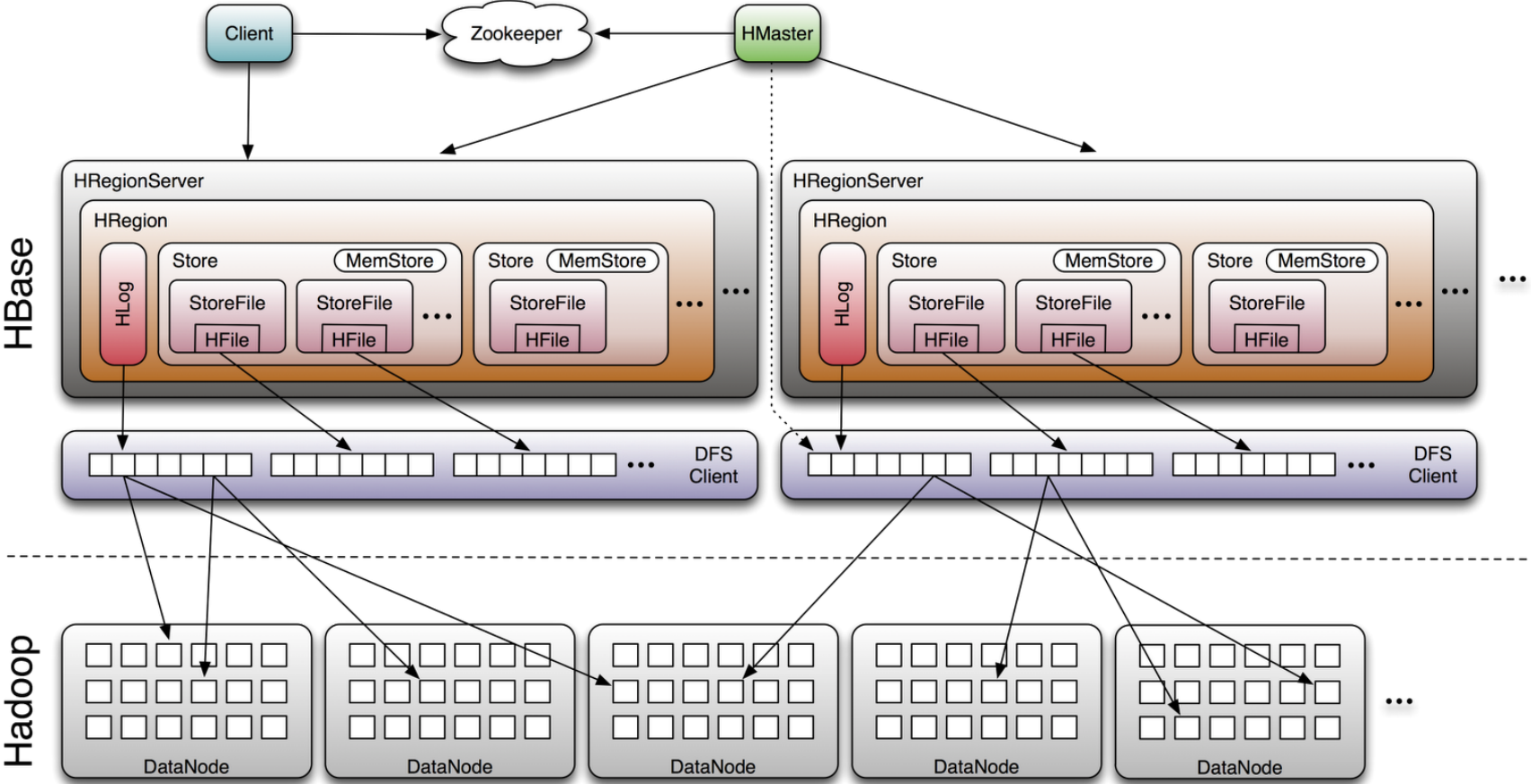
The Department Administrator



WWW.PHDCOMICS.COM

Morale of the story: there's no free lunch!

HBase



Three Core Ideas

- Partitioning (sharding)

- For scalability
- For latency

Need distributed transactions!

- Replication

- For robustness (availability)
- For throughput

Need replica coherence protocol!

- Caching

- For latency

Need cache coherence protocol!

This is really hard!



Now imagine multiple datacenters...
What's different?

Three Core Ideas

- Partitioning (sharding)

- For scalability
- For latency

Need distributed transactions!

- Replication

- For robustness (availability)
- For throughput

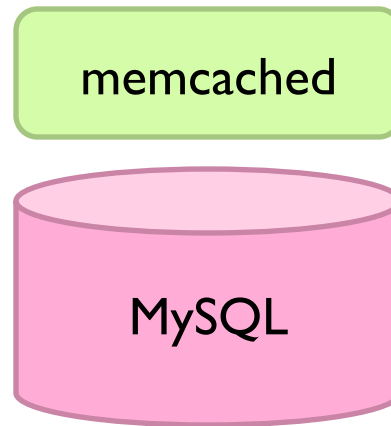
Need replica coherence protocol!

- Caching

- For latency

Need cache coherence protocol!

Facebook Architecture



Read path:

Look in memcached
Look in MySQL
Populate in memcached

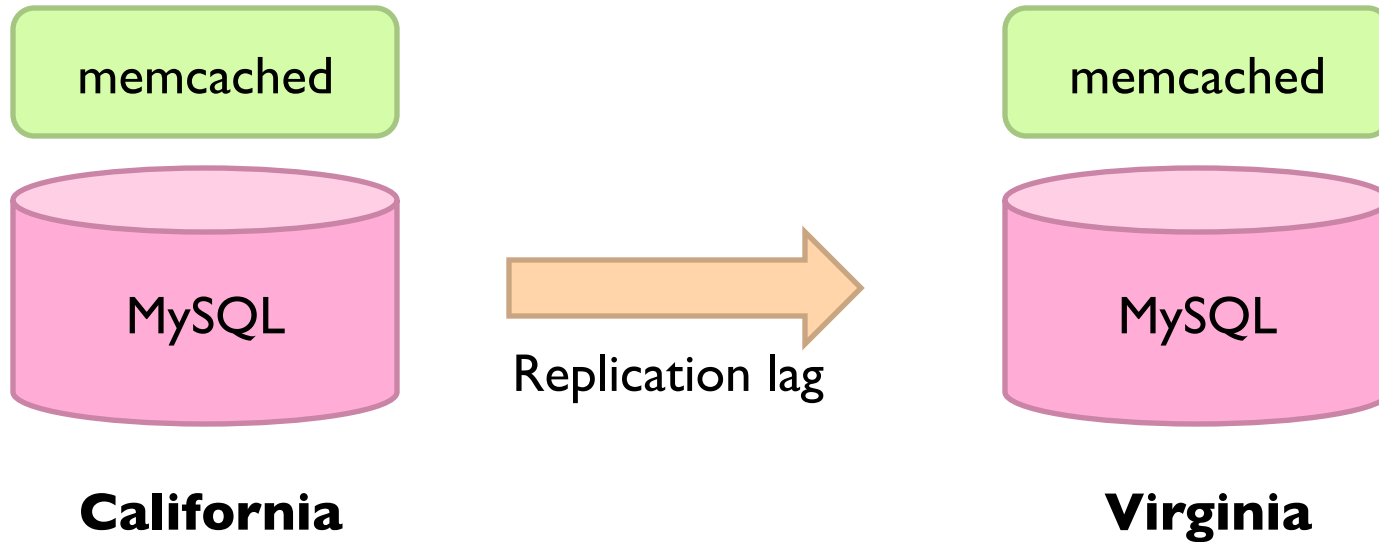
Write path:

Write in MySQL
Remove in memcached

Subsequent read:

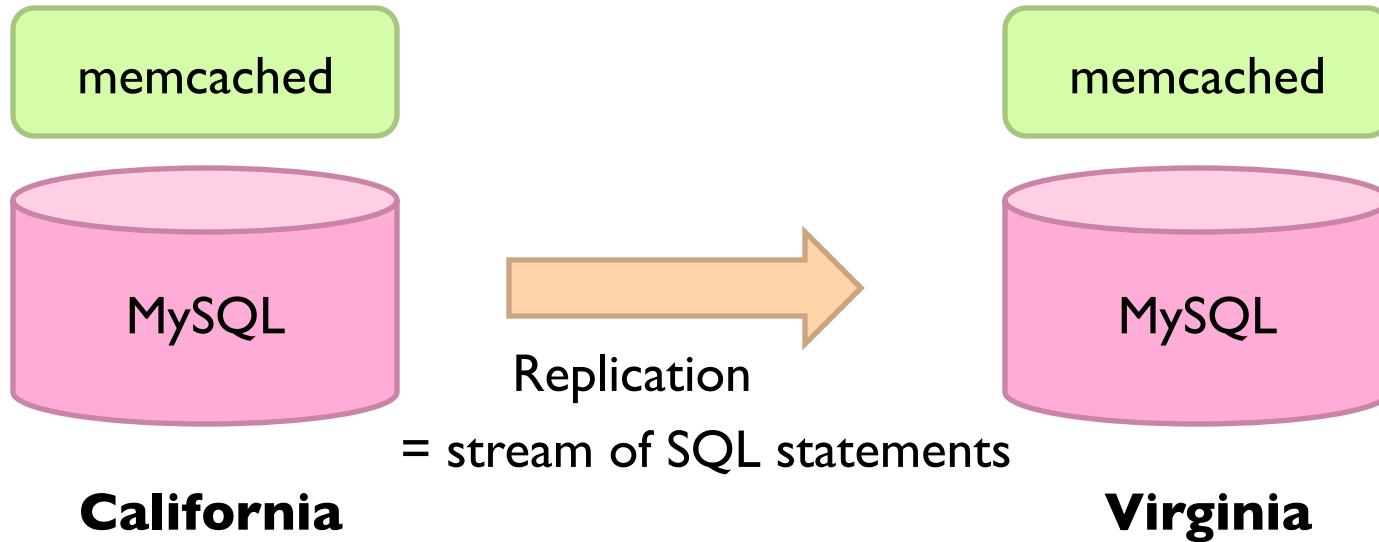
Look in MySQL ✓
Populate in memcached

Facebook Architecture: Multi-DC



1. User updates first name from “Jason” to “Monkey”.
2. Write “Monkey” in master DB in CA, delete memcached entry in CA and VA.
3. Someone goes to profile in Virginia, read VA slave DB, get “Jason”.
4. Update VA memcache with first name as “Jason”.
5. Replication catches up. “Jason” stuck in memcached until another write!

Facebook Architecture



Solution: Piggyback on replication stream, tweak SQL

```
REPLACE INTO profile (`first_name`) VALUES ('Monkey')  
WHERE `user_id`='jsobel' MEMCACHE_DIRTY 'jsobel:first_name'
```

Three Core Ideas

- Partitioning (sharding)

- For scalability
- For latency

Need distributed transactions!

- Replication

- For robustness (availability)
- For throughput

Need replica coherence protocol!

- Caching

- For latency

Need cache coherence protocol!

Yahoo's PNUTS

- Yahoo's globally distributed/replicated key-value store
- Provides *per-record* timeline consistency
 - Guarantees that all replicas provide all updates in same order
- Different classes of reads:
 - Read-any: may time travel!
 - Read-critical(required version): monotonic reads
 - Read-latest

PNUTS: Implementation Principles

- Each record has a single master
 - Asynchronous replication across datacenters
 - Allow for synchronous replicate within datacenters
 - All updates routed to master first, updates applied, then propagated
 - Protocols for recognizing master failure and load balancing
- Tradeoffs:
 - Different types of reads have different latencies
 - Availability compromised when master fails and partition failure in protocol for transferring of mastership

Three Core Ideas

○ Partitioning (sharding)

- For scalability
- For latency

Need distributed transactions!

○ Replication

- For robustness (availability)
- For throughput

Have our cake and eat it too?

Need replica coherence protocol!

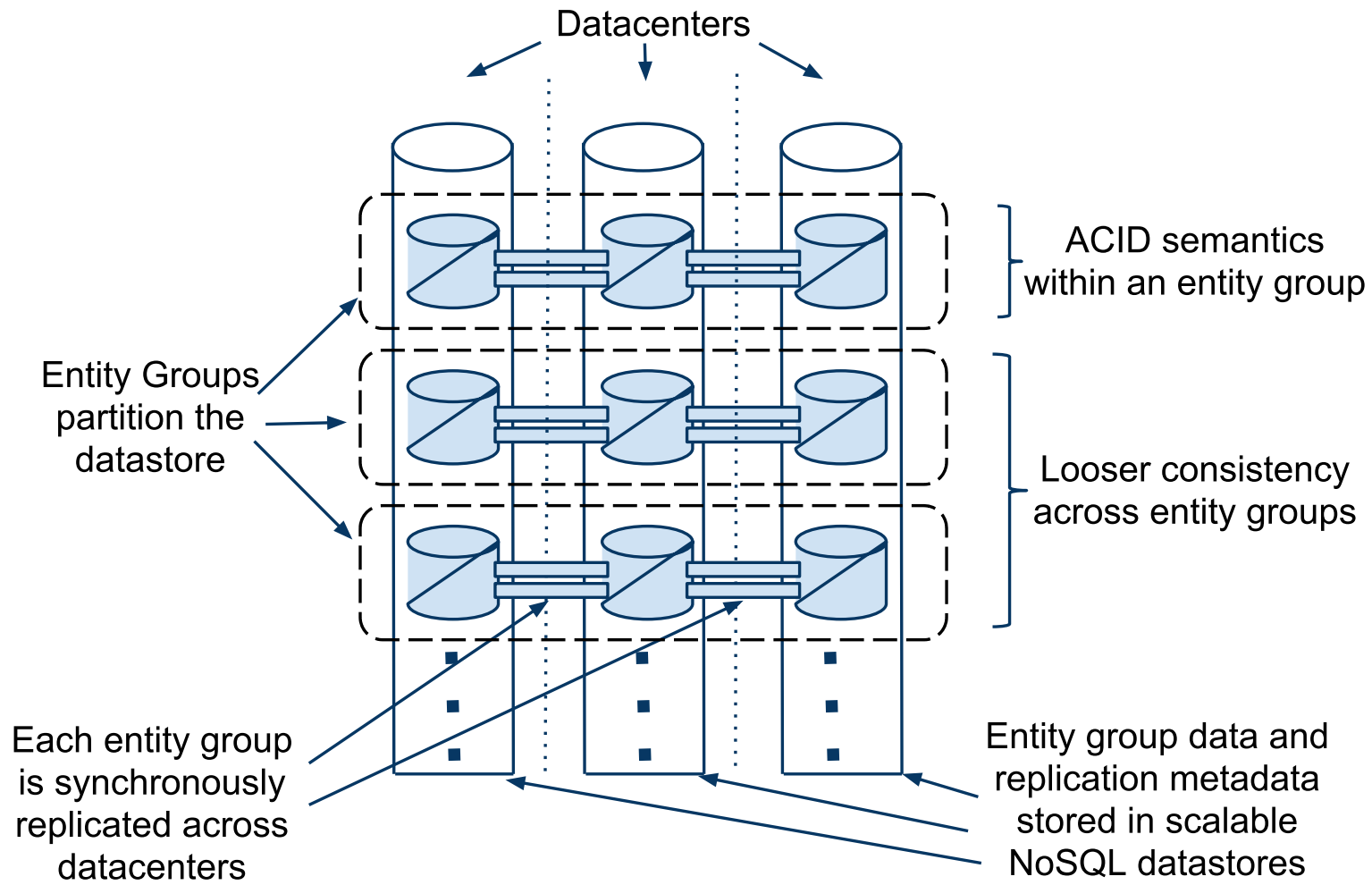
○ Caching

- For latency

Need cache coherence protocol!



Google's Megastore



Google's Spanner

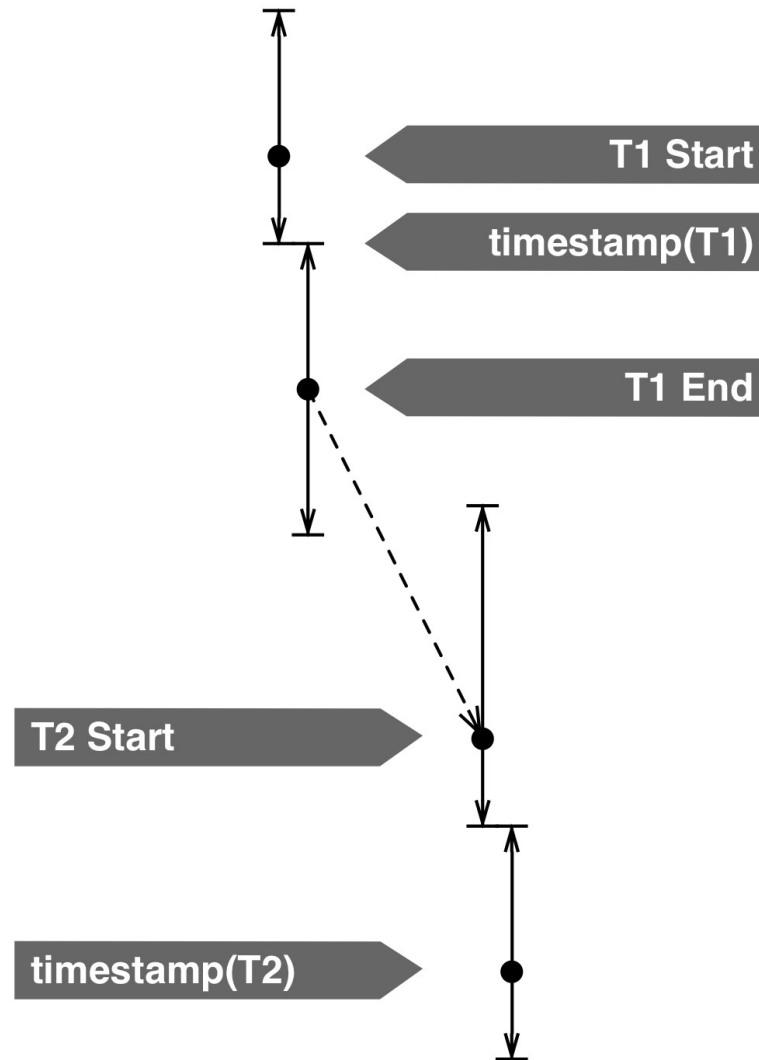
○ Features:

- Full ACID translations across multiple datacenters, across continents!
- External consistency (= linearizability):
system preserves *happens-before* relationship among transactions

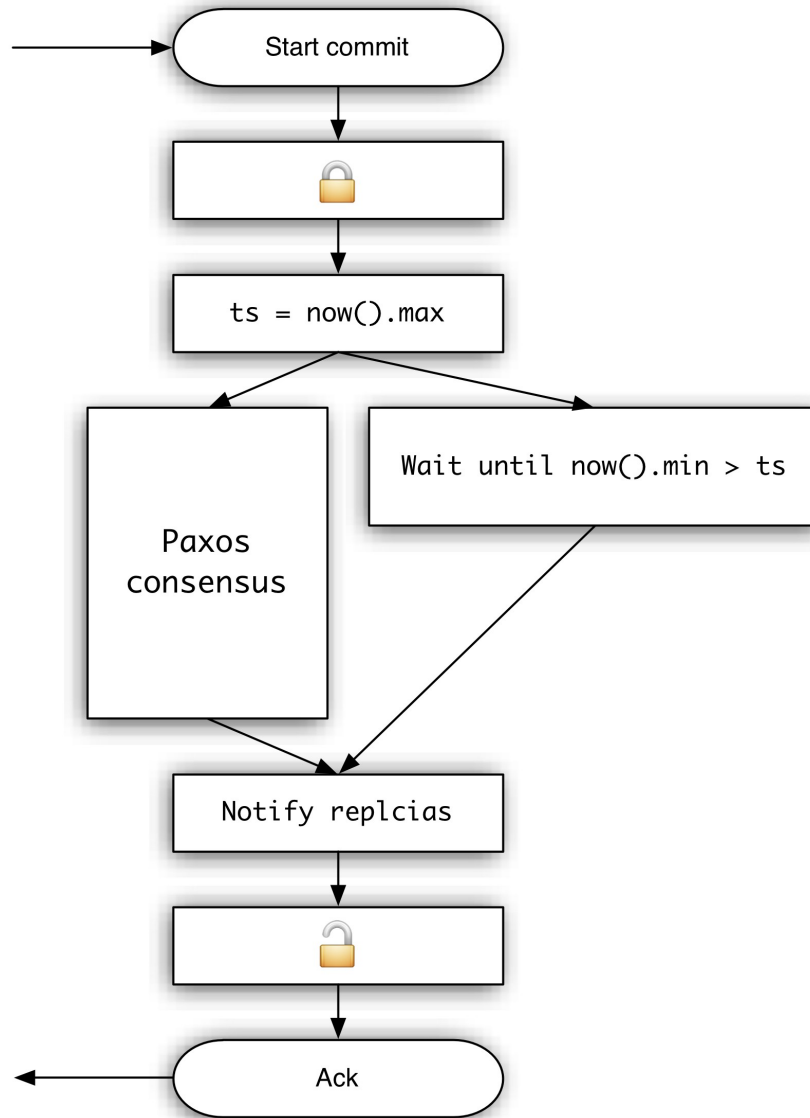
○ How?

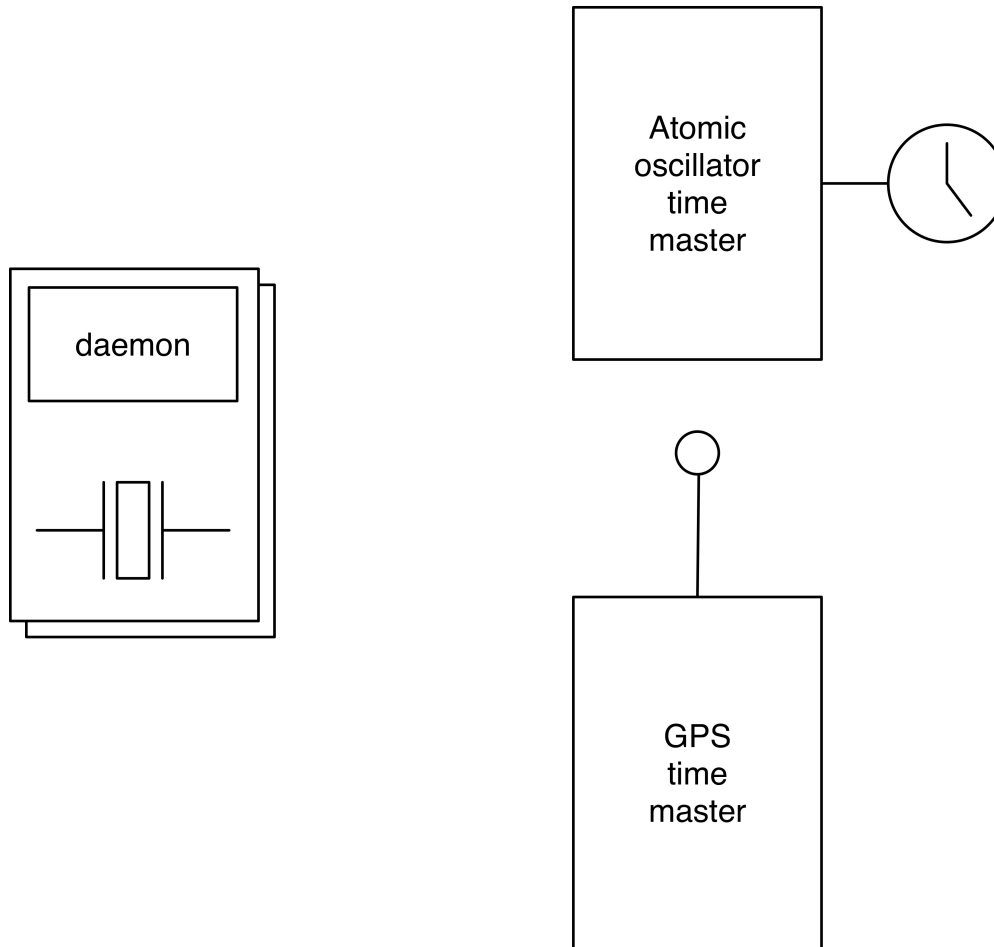
- Given write transactions A and B, if A *happens-before* B, then $\text{timestamp}(A) < \text{timestamp}(B)$

Why this works



TrueTime → write timestamps







What's the catch?

To: All Graduate Students

Due to a recent incident, we would like to remind all Grad Students that refreshments provided in communal areas during an event are for attendees of that event only.

Please vacate the communal area and do not consume the refreshments unless you have been specifically invited to participate.

To avoid any misunderstanding, you are only invited if you received a specific invitation by e-mail or if it was arranged by your supervisor for you to attend.

Thank you for your cooperation,

The Department Administrator



WWW.PHDCOMICS.COM

Morale of the story: there's no free lunch!



Questions?